

OCR von Inkunabeln: Herausforderungen und Herangehensweisen

Uwe Springmann

Centrum für Informations- und Sprachverarbeitung (CIS)
Ludwig-Maximilians-Universität München
und
Institut für deutsche Sprache und Literatur
Humboldt-Universität zu Berlin

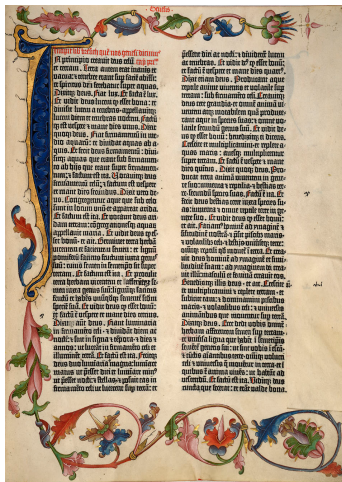


<philtag n="13"/>, Universität Würzburg

2016-02-25

Einleitung

Warum interessieren uns Inkunabeln (Wiegendrucke)?



als die Druckkunst in der Wiege lag
(lat. incunabula, orum n.: Wiege)

- älteste Dokumente der modernen Druckgeschichte (1450-1500)
- Medienrevolution fällt in Umbruchzeit:
 - Vorabend der Reformation (Luther 1517)
 - Entdeckung Amerikas (Kolumbus 1492)
 - Wiederentdeckung der Antike (Renaissance)

Gutenberg-Bibel B42, 1454

Inkunabeln sind spannende Quellen einer
Schlüsselperiode Europas!

Herausforderungen

Gestaltung und Typographie



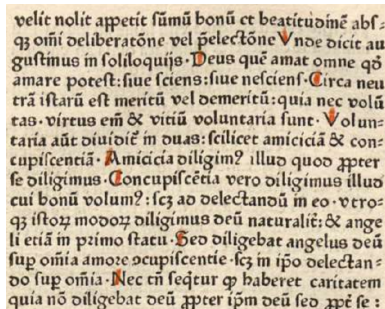
Rothschild-Bibel, 13. Jh.,
Paris

- absatzmarktorientierter Anschluss an Gestaltung von Manuskripten (illuminiertes Manuskript kostete soviel wie ein Haus)
- Farbe, Zierinitialen, Schmuck, Ligaturen, Sonderzeichen
- Papierbewirtschaftung: Spaltendruck, zahlreiche Abkürzungen
- noch keine Arbeitsteilung und Standardisierung: jeder Frühdrucker war auch Typenschneider und Setzer (Typenvielfalt!)

Hürden für automatische Segmentierung und Transkription!

OCR auf Inkunabeln: status quæstionis (bis 2014))

Beauvais: *Speculum naturale* (nicht nach 1476); ABBYY FR11 Fraktur 68% acc.



velit nolit appetit sumū bonū et beatitudinē abs-
q3 omī deliberatōne vel p̄lectōne Vnde dicit au-
gustinus in soliloquijs. Deus quē amat omne q3
amare potest: siue sciens: siue nesciens. Circa neu-
trā istarū est meritū vel demeritū: quia nec volū-
tas. virtus enī & viciū voluntaria sunt. Volun-
taria autē diuidit̄ in duas: scilicet amicitia & con-
cupiscenciā. Amicitia diligim? illud quod ppter
se diligimus. Concupiscenciā vero diligimus illud
cui bonū volum?: sc3 ad delectandū in eo. vero-
q3 istoꝝ modoz diligimus deū naturalit̄: & ange-
li etiā in primo statu. Sed diligebat angelus deū
sup omīa amore & cupiscencie. sc3 in ip̄o delectan-
do sup omīa. Nec tñ sequitur q3 haberet caritatem
quia nō diligebat deū ppter ip̄m deū sed ppter se:

velie nolit aspenc sumu bonu ce beatituome al? s-
qzonn veliberaeone velpelec^oneVnoeivicau
Aus^mus in soliloquijs'^eus que amat omne qv
amarc potest:s»uesciens.smenesciens Circa neu
era ls^aru esk mencu vet vcmenturquia ncc volu
ras vireus em viciu voluntana (une V^o!un-
tana auc oiuivic in ouas: scilicet amicitia Le con-
cupilcencia 5Vmicicra vilizim? illuo quov zpter
sevoli^imus Concupiscena vcro viliquimus illuo
cui bonu volum?:lc3 as velec^anvu in co Vtro-
qz is^or^ movoy oiliquimus veu naeuralitrLL an^e
li eria in primis l^acu Leo viliquebat an^clus veu
sup omia amorc occupiscentie lc3 in ipo vclec^an -
vo sup omia Alec cn sequer q? kaberee caritacem
quia no viligvbat veu ^fpter ipm veu seo zx>c se :

Inkunabeln haben häufig besondere Abkürzungszeichen, z.B. p p p q Q q sc3.

(Rydberg-Cox 2009) (unsere Hervorhebung): “Because of the prevalence of these glyphs, *incunabula cannot be processed using OCR software*. Commercial OCR programs produce almost no recognizable character strings, let alone searchable text. ... *Other methods must be explored.*”

Herangehensweisen

Andere (OCR-) Methoden: Rekurrente neuronale Netze

- Schlagwort:
Rekurrente neuronale Netze mit langem Kurzzeitgedächtnis
RNN mit LSTM, *Hochreiter and Schmidhuber (1997)*
- Methode hatte große Erfolge bei Mustererkennung
(Verkehrszeichen, Gesichtserkennung, ...)
- auf OCR-Erkennung erstmals angewandt von *Breuel et al. (2013)*
- auf Erkennung von Frühdrucken adaptiert von *Springmann et al. (2014)*;
Springmann, Lüdeling, and Schremmer (2015); *Springmann (2015)*

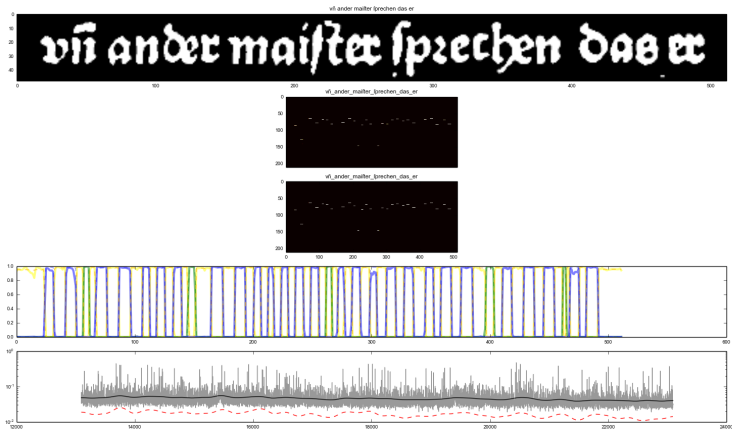
Wie lernt das neuronale Netz?

Idee (Breuel):

- zerschneide Bild einer Textzeile in viele vertikale Streifen (500-1000)
- ordne den Streifen (Pixels) einer Zeile die diplomatische Transkription (Labels) der Zeile zu
- das Netzwerk gewichtet die Verbindungen seiner internen Speicherzellen (Gedächtnis) so, dass eine Verbindung von Inputdaten (Pixels) zu Outputdaten (Labels) entsteht
- Lernen geschieht selbsttätig (Klassifizieren von benachbarten Streifen zu kodierten Glyphen)
- nach einiger Zeit erkennt es vorher nicht gesehene Zeilen mit guter Genauigkeit

Die Zerlegung von Zeichen in einzelne Streifen als Grundeinheiten ist der Schlüssel für die bessere Erkennung gegenüber einer Mustererkennung auf Zeichenebene!

Dem Netz beim Lernen zuschauen



Trainieren eines OCR-Modells auf einer Schriftart

Das Modelltraining gliedert sich in die folgenden Schritte:

- ❶ Beschaffen der Scans
- ❷ Zerlegen der Seitenbilder in einzelne Zeilen
- ❸ Herstellung einer diplomatischen Transkription (*ground truth*) dieser Zeilen
- ❹ Aufteilen der Bild- und Textzeilen in eine Trainings- und eine Testmenge
- ❺ Training auf der Trainingsmenge
- ❻ Testen auf Testmenge:
 - Testergebnis ok: Erkennen des ganzen Dokumentes
 - Testergebnis zu schlecht:
Korrektur der Erkennung einiger weiterer Seiten zu ground-truth-Qualität,
Hinzufügen zur Trainingsmenge und Rücksprung auf Nr. 5

Diplomatische Transkription: *ground truth*-Erstellung

Beyfuß

Beyfuß

Das Erft Capitel

Das Erft Capitel

Ariuofa Ampolata Brita

Ariuofa Ampolata Brita

nica Campanaria Metri

nica Campanaria Metri

caria minor: lat̃e·Melenoff Zans

caria minor lat̃e + Melenoff Zans

tes Thagetes Leptafelos ¶ Die

tes Thagetes Leptafelos ¶ Die

wirdigen maister Auicenna Dia

wirdigen maister Auicenna Dia

scorides beschreiben vns vō difem

scorides beschreiben vns vō difem

- Eingabe über Zeilensynopse im Browser mit geeigneter Schriftart, z.B. [Junicode](#)
- Glyph-Repertoire bestimmen ([Häberle](#))
- Paläographie-Kenntnisse notwendig
 - Ligaturen
 - Suspensionen
 - Kontraktionen
- weitere Voraussetzungen:
 - historische Linguistik
 - Schreibvarianten
 - Frühneuhochdeutsch
 - Latein ([70% der Inkunabeldrucke](#))

Modelltraining

nach einer Weile (hier: nach 23.055 Lernschritten):

23055 1.99 (497, 48) **train/0004/01001f.bin.png**

TRU: u'vertreibt die \u017fchlangen die in den'

ALN: u'vertreibt die \u017fchlangen die in den'

OUT: u'vertreibt die \u017fchlan gen die in den'

23056 1.42 (508, 48) **train/0002/010046.bin.png**

TRU: u'laxieren i\u017ft vnd purgieren / das dz'

ALN: u'laxieren i\u017ft vnd purgieren / das dz'

OUT: u'laxieren i\u017ft vnd purgieren / das dz'

23057 2.02 (514, 48) **train/0001/01002e.bin.png**

TRU: u'che fraw wee mit ainem kind gat'

ALN: u'che fraw wee mit ainem kind gat'

OUT: u'che fra w wee mit ainem kind gat'

Noch einmal Beauvais, *Speculum Naturale*

Trainiertes OCRopus-Modell (dieser Ausschnitt: 99% acc.)

velit nolit appetit sūmū bonū et beatitudinē abf-
q3 omī deliberatōne vel p̄lectōne Vnde dicit au-
gustinus in soliloquijs · Deus quē amat omne qđ
amare potest: siue sciens: siue nesciens · Circa neu-
trā istarū est meritū vel demeritū: quia nec volū-
tas · virtus em̄ & vitiū voluntaria sunt · Volun-
taria aut̄ diuidit̄ in duas: scilicet amiciciā & con-
cupiscenciā · Amicicia diligim⁹ illud quod ppter
se diligimus · Concupiscenciā vero diligimus illud
cui bonū volum⁹: sc3 ad delectandū in eo · vtro-
q3 istoꝝ modoz diligimus deū naturalit̄: & ange-
li etiā in primo statu · Sed diligebat angelus deū
sup omīa amore concupiscentie · sc3 in ip̄o delectan-
do sup omīa · Nec tñ seq̄tur qđ haberet caritatem
quia nō diligebat deū ppter ip̄m deū sed ppter se :

velit nolit appetit sūmū bonū et beatitudinē abf-
q3 omī deliberatōne vel p̄lectōne Vnde dicit au-
gustinus in foliloquijs · Deus quē amat omne qđ
amare potest: siue sciens: siue nesciens · Circa neu-
trā istarū est meritū vel demeritū : quia nec volū-
tas · virtus em̄ & vitiū voluntaria sunt · Volūne-
taria aut̄ diuidit̄ in duas: scilicet amicicia & con-
cupiscenciā · Amicicia diligim⁹ illud quod ppter
se diligimus · Concupiscenciā vero diligimus illud
cui bonū volum⁹ : sc3 ad delectandū in eo · vtro-
q3 istoꝝ modoꝝ diligimus deū naturalit̄: & ange-
li etiā in primo statu · Sed diligebat angelus deū
sup omīa amore concupiscentie · f3 in ip̄o delectan-
do sup omīa · Nec tñ seq̄tur qđ haberet caritatem
quia nō diligebat deū ppter ip̄m deū sed ppter se :

- nur noch 4 Fehler! (rechts: rot und blau markiert)
- trainiert auf 13 Seiten, getestet auf weiteren 4 Seiten
- 98% mittlere Zeichenerkennungsrate (rohes, unkorrigiertes OCR-Ergebnis)
- ohne Verwendung eines Sprachmodells

Offene Fragen

Muss man jede Type separat trainieren?

	1487-Gar	1532-Art	1532-Cor	1543-Neu	1557-Wie	1588-Par	1603-Alc	1609-Kra	1639-Pfla	1652-Wui	1673-The	1675-Cur	1687-Der	1735-My	1764-Einl	1774-Unt	1828-Die	1870-Deu
1487-GartDe	3,522	27,513	29,79	24,506	18,514	27,384	26,632	28,78	23,024	28,351	28,501	24,871	32,71	28,565	33,57	28,479	37,264	48,282
1532-Artzney	17,594	1,109	16,464	13,096	10,335	11,444	9,372	13,452	11,925	12,05	15,879	10,921	17,678	19,937	24,916	24,665	31,318	38,368
1532-Contraf	26,709	27,32	3,671	18,388	18,056	21,727	19,175	27,338	23,72	35,763	35,011	23,422	46,967	29,628	33,298	34,382	46,128	53,994
1543-NewKre	12,091	12,484	13,644	1,853	10,313	14,542	12,952	15,815	12,596	17,705	19,053	16,732	20,251	17,219	25,08	26,745	33,502	37,226
1557-WieSci	13,247	20,786	16,618	14,68	9,003	14,362	17,813	23,334	14,176	17,069	24,821	24,821	27,582	21,184	24,078	19,166	27,874	42,952
1588-Paradei	32,801	35,505	32,698	36,071	29,016	2,652	14,779	17,173	10,736	10,273	12,796	10,891	16,89	20,881	26,571	23,532	29,763	37,204
1603-Alchym	25,595	22,794	20,53	26,551	24,925	10,013	2,986	8,017	7,43	13,418	12,965	11,338	18,719	14,24	25,948	21,419	38,007	33,814
1609-Kraeut	19,321	22,935	17,503	20,504	16,583	5,476	3,023	1,402	6,221	7,514	7,426	6,55	9,398	12,574	18,905	18,335	26,287	27,952
1639-Pflantz	21,093	23,525	23,426	24,808	19,01	8,181	7,348	10,113	2,416	4,898	9,563	6,548	13,512	13,112	18,077	13,545	24,442	28,607
1652-Wund-A	30,476	26,692	32,772	44,374	24,483	8,573	13,756	16,173	5,631	2,515	10,202	8,398	12,761	16,271	21,432	17,321	26,9	32,881
1673-Thesau	30,476	25,833	29,868	32,096	27,118	8,711	10,739	10,458	8,41	7,367	1,044	6,042	9,494	11,642	14,151	17,623	23,906	26,716
1675-Curiose	30,476	29,105	28,148	25,741	23,21	10,988	16,265	17,037	7,623	10,278	14,63	5	21,605	20	22,377	19,29	26,204	32,006
1687-DerSch	30,476	30,188	29,662	43,571	34,436	16,617	14,135	13,459	12,068	10,15	7,481	5,414	3,233	15,977	15,188	17,82	24,248	29,135
1735-Mysteri	30,476	24,794	21,765	23,214	23,872	12,315	12,94	15,608	7,935	10,998	9,121	8,133	13,006	2,14	9,055	8,331	15,838	15,904
1764-Einleit	30,476	27,075	24,127	29,428	23,869	7,422	12,982	13,24	6,956	7,448	5,017	2,974	8,379	5,094	0,362	3,698	6,206	10,654
1774-Unterri	30,476	24,019	17,813	24,145	15,829	6,754	7,767	10,426	5,572	6,036	9,329	3,926	9,118	4,052	4,137	0,549	6,079	14,141
1828-DieEige	30,476	28,773	22,458	31,84	23,063	13,943	18,805	21,57	11,481	11,461	13,963	6,174	15,133	8,717	6,174	7,425	0,605	8,535
1870-Deutsci	30,476	35,064	29,749	41,771	33,245	18,299	22,429	21,407	15,907	16,152	14,803	11,429	17,134	15,457	12,86	13,883	13,126	1,738

Frakturmodelle angewendet auf Frakturdrucke des **RIDGES-Korpus**. Die Zahlen geben den Fehler der OCR-Erkennung auf Zeichenebene an (Prozentsatz falsch erkannter Zeichen).

- Spalten: Modelle
- Zeilen: Drucke
- beste Ergebnisse auf Diagonale (Modell passt zum Druck)
- Nichtdiagonalelemente: bei Frühdrucken starke Variation, später weniger
- keine systematische Untersuchung, durch verfügbares Material begrenzt
- **Derzeit mangels Modellvielfalt / Trainingsdatenknappheit noch keine klare Aussage möglich!**

Jenseits der OCR-Erkennung: Normalisierung

- Wie sucht man auf diesem Text?

Von di

sem kraut beschreibet vns Diasco-
rides vnd sprichet + das dñes kraut
beneme vnnð haile acrocordines
das sind lychdorn oð wärçzẽ auff

- Normalisierung auf heutige Schreibweise als Annotationsbeene notwendig, siehe dazu z.B. Bollmann, Petran, and Dipper (2011), Jurish (2013)

Jenseits der OCR-Erkennung: Nachkorrektur

Nachkorrektur: z.B. mit dem interaktiven CIS-Tool **PoCoTo**

The screenshot displays the PoCoTo (Pocoto) software interface, designed for the interactive correction of OCR errors in historical documents. The main window shows a page of Greek text from the New Testament (John 1:1-14) with various tokens highlighted for correction. The interface includes a 'Concordance Actions' sidebar on the left with buttons for 'Show concordance', 'Multi token actions', 'Merge selected tokens', and 'Delete selected tokens'. The main text area is divided into columns, each showing a token, its original OCR, and a suggested correction. For example, 'τοῖς δικαίοις καὶ ἀγαθοῖς' is corrected to 'τοῖς δικαίοις καὶ ἀγαθοῖς'. The bottom of the window shows the original Greek text with the corrected tokens integrated. The interface also includes a 'Page 1 of 846' indicator and a search bar in the top right corner.

Vorschläge für ein koordiniertes Vorgehen

- Verfügbarmachung der vorhandenen Scans von Buchseiten in hochaufgelöster Form (nicht binarisierte tiffs anstelle binarisierter jpgs) für die Wissenschaft
- Akquisition eines Budgets zur Aufwertung der vorhanden Scanbestände durch OCR (für einen Bruchteil des für Scans ausgegebenen Budgets gewinnt man ein Mehrfaches an Nachnutzungspotential)
- Einrichtung eines gemeinsamen Daten-Repositories für OCR-Daten, die mit anderen bestandshaltenden Institutionen geteilt und zentral über eine Webschnittstelle zur Nachkorrektur angeboten werden (auch für *crowd-sourcing* geeignet)
- Zentrales Modelltraining in einem leistungsfähigen Rechenzentrum sowohl auf Einzeltypographien als auch für Typenmischungen, um die Einsetzbarkeit der Modelle zu verbessern
- Koordiniertes Vorgehen anstatt jeder für sich dasselbe machen (wie bei Scans)
- Ein solches Projekt würde eine weltweit einmalige Forschungsgrundlage auf Basis der Drucke aus deutschen Gebieten darstellen (einschließlich der in deutschen Gebieten gedruckten lateinischen etc. Werke)

Fazit

- Wir können heute bereits Drucke der gesamten modernen Druckgeschichte bis hinunter zu Gutenberg mit hoher Genauigkeit (> 95%) durch eine typentrainierte OCR erkennen.
- Es fehlt an verlässlicher *ground truth* und an einer Kultur des offenen Datenaustausches. Auch die Urheberrechtsfrage von Scans (Stichwort *copyfraud*, *Schutzrechtsberühmung* bei gemeinfreien Inhalten) und die damit einhergehende Nichtverfügbarkeit hochauflöster Bilddaten erschwert das Modelltraining.
- Eine koordinierte Initiative deutscher Institutionen könnte auf Basis der vorliegenden Scans eine OCR-Erfassung durchführen und zentral zur Nachkorrektur anbieten. Der dadurch entstehende *ground truth*-Vorrat könnte automatisiert zur Modellverbesserung genutzt und damit ein sich stets verbessernder Zirkel in Gang gesetzt werden, an dessen Ende das gesamte bildmäßig erfasste Material als hochgenauer elektronischer Text vorliegt.

Wenn Sie mehr erfahren möchten

- [CIS OCR Workshop](#) (Springmann and Fink 2016)
- [Ocrocis](#) (Springmann and Kaumanns 2015)
A project manager interface to OCRopus
- [Ocrocis Tutorial](#) (Springmann 2015)
Ausführliche Anleitung zum Trainieren eigener Modelle
- ein allgemeines [OCR Tutorial](#) (Springmann 2014)

Vielen Dank für Ihre Aufmerksamkeit!

Dr. Uwe Springmann
§ digital humanist §
vorname [A T] nachname.net

Literaturangaben I

Bollmann, Marcel, Florian Petran, and Stefanie Dipper. 2011. “Applying Rule-Based Normalization to Different Types of Historical Texts—an Evaluation.” In *Human Language Technology Challenges for Computer Science and Linguistics*, 166–77. Springer.

Breuel, Thomas M, Adnan Ul-Hasan, Mayce Ali Al-Azawi, and Faisal Shafait. 2013. “High-Performance OCR for Printed English and Fraktur Using LSTM Networks.” In *2th International Conference on Document Analysis and Recognition (ICDAR), 2013*, 683–87. IEEE.

Hochreiter, Sepp, and Jürgen Schmidhuber. 1997. “Long Short-Term Memory.” *Neural Computation* 9 (8). MIT Press: 1735–80.

Jurish, Bryan. 2013. “Canonicalizing the Deutsches Textarchiv.” In *Proceedings of Perspektiven Einer Corpusbasierten Historischen Linguistik Und Philologie (Berlin, 12th - 13th December 2011)*, edited by Ingelore Hafemann. Vol. 4. Thesaurus Linguae

Literaturangaben II

Aegyptiae. Berlin, Germany: Berlin-Brandenburgische Akademie der Wissenschaften.
http://edoc.bbaw.de/frontdoor.php?source_opus=2443.

Rydberg-Cox, Jeffrey A. 2009. "Digitizing Latin Incunabula: Challenges, Methods, and Possibilities." *Digital Humanities Quarterly* 3 (1).
<http://www.digitalhumanities.org/dhq/vol/3/1/000027/000027.html/#p7>.

Springmann, Uwe. 2015. "Ocrocis: A high accuracy OCR method to convert early printings into digital text – A Tutorial." <http://cistern.cis.lmu.de/ocrocis/tutorial.pdf>.

Springmann, Uwe, and Florian Fink. 2016. "CIS OCR Workshop v1.0: OCR and postcorrection of early printings for digital humanities." [doi:10.5281/zenodo.46571](https://doi.org/10.5281/zenodo.46571).

Springmann, Uwe, and David Kaumanns. 2015. "Ocrocis – a high accuracy OCR method to convert early printings into digital text." <http://cistern.cis.lmu.de/ocrocis/>.

Literaturangaben III

Springmann, Uwe, Anke Lüdeling, and Felix Schremmer. 2015. “Zur OCR frühneuzeitlicher Drucke am Beispiel des RIDGES-Korpus von Kräutertexten.” DHd-Tagung 2015, Graz. <http://gams.uni-graz.at/o:dhd2015.p.34>.

Springmann, Uwe, Dietmar Najock, Hermann Morgenroth, Helmut Schmid, Annette Gotscharek, and Florian Fink. 2014. “OCR of historical printings of Latin texts: problems, prospects, progress.” In *Proceedings of the First International Conference on Digital Access to Textual Cultural Heritage*, 57–61. DATeCH '14. New York, NY, USA: ACM. [doi:10.1145/2595188.2595197](https://doi.org/10.1145/2595188.2595197).